

REPORT DOCUMENTATION PAGE			Form Approved OMB NO. 0704-0188		
<p>The public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA, 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to any penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.</p> <p>PLEASE DO NOT RETURN YOUR FORM TO THE ABOVE ADDRESS.</p>					
1. REPORT DATE (DD-MM-YYYY)		2. REPORT TYPE New Reprint		3. DATES COVERED (From - To) -	
4. TITLE AND SUBTITLE Structure, function and diversity of the healthy human microbiome			5a. CONTRACT NUMBER W911NF-11-1-0429		
			5b. GRANT NUMBER		
			5c. PROGRAM ELEMENT NUMBER 611102		
6. AUTHORS Krishna Palaniappan, Shital M. Patel, Matthew Pearson, Jane Peterson, Mircea Podar, Craig Pohl, Katherine S. Pollard, Doyle V. Ward, Wesley Warren, Mark A. Watson, Christopher Wellington, Kris A.			5d. PROJECT NUMBER		
			5e. TASK NUMBER		
			5f. WORK UNIT NUMBER		
7. PERFORMING ORGANIZATION NAMES AND ADDRESSES Harvard School of Public Health Biostatistics President and Fellows of Harvard College Boston, MA 02115 -6028				8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES) U.S. Army Research Office P.O. Box 12211 Research Triangle Park, NC 27709-2211				10. SPONSOR/MONITOR'S ACRONYM(S) ARO	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S) 60119-MA.3	
12. DISTRIBUTION AVAILABILITY STATEMENT Approved for public release; distribution is unlimited.					
13. SUPPLEMENTARY NOTES The views, opinions and/or findings contained in this report are those of the author(s) and should not be construed as an official Department of the Army position, policy or decision, unless so designated by other documentation.					
14. ABSTRACT Studies of the human microbiome have revealed that even healthy individuals differ remarkably in the microbes that occupy habitats such as the gut, skin and vagina. Much of this diversity remains unexplained, although diet, environment, host genetics and early microbial exposure have all been implicated. Accordingly, to characterize the ecology of human-associated microbial communities, the Human Microbiome Project has analysed the largest cohort and set of distinct, clinically relevant body habitats so far. We found the diversity and abundance of each					
15. SUBJECT TERMS human microbiome					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT UU	15. NUMBER OF PAGES	19a. NAME OF RESPONSIBLE PERSON Curtis Huttenhower
a. REPORT UU	b. ABSTRACT UU	c. THIS PAGE UU			19b. TELEPHONE NUMBER 617-432-4912

Report Title

Structure, function and diversity of the healthy human microbiome

ABSTRACT

Studies of the human microbiome have revealed that even healthy individuals differ remarkably in the microbes that occupy habitats such as the gut, skin and vagina. Much of this diversity remains unexplained, although diet, environment, host genetics and early microbial exposure have all been implicated. Accordingly, to characterize the ecology of human-associated microbial communities, the Human Microbiome Project has analysed the largest cohort and set of distinct, clinically relevant body habitats so far. We found the diversity and abundance of each habitat's signature microbes to vary widely even among healthy subjects, with strong niche specialization both within and among individuals. The project encountered an estimated 81–99% of the genera, enzyme families and community configurations occupied by the healthy Western microbiome. Metagenomic carriage of metabolic pathways was stable among individuals despite variation in community structure, and ethnic/racial background proved to be one of the strongest associations of both pathways and microbes with clinical metadata. These results thus delineate the range of structural and functional configurations normal in the microbial communities of a healthy population, enabling future characterization of the epidemiology, ecology and translational applications of the human microbiome.

REPORT DOCUMENTATION PAGE (SF298)
(Continuation Sheet)

Continuation for Block 13

ARO Report Number 60119.3-MA

Structure, function and diversity of the healthy hu ...

Block 13: Supplementary Note

© 2012 . Published in Nature, Vol. Ed. 0 486, (7402) (2012), (7402). DoD Components reserve a royalty-free, nonexclusive and irrevocable right to reproduce, publish, or otherwise use the work for Federal purposes, and to authroize others to do so (DODGARS §32.36). The views, opinions and/or findings contained in this report are those of the author(s) and should not be construed as an official Department of the Army position, policy or decision, unless so designated by other documentation.

Approved for public release; distribution is unlimited.

Structure, function and diversity of the healthy human microbiome

The Human Microbiome Project Consortium*

Studies of the human microbiome have revealed that even healthy individuals differ remarkably in the microbes that occupy habitats such as the gut, skin and vagina. Much of this diversity remains unexplained, although diet, environment, host genetics and early microbial exposure have all been implicated. Accordingly, to characterize the ecology of human-associated microbial communities, the Human Microbiome Project has analysed the largest cohort and set of distinct, clinically relevant body habitats so far. We found the diversity and abundance of each habitat's signature microbes to vary widely even among healthy subjects, with strong niche specialization both within and among individuals. The project encountered an estimated 81–99% of the genera, enzyme families and community configurations occupied by the healthy Western microbiome. Metagenomic carriage of metabolic pathways was stable among individuals despite variation in community structure, and ethnic/racial background proved to be one of the strongest associations of both pathways and microbes with clinical metadata. These results thus delineate the range of structural and functional configurations normal in the microbial communities of a healthy population, enabling future characterization of the epidemiology, ecology and translational applications of the human microbiome.

A total of 4,788 specimens from 242 screened and phenotyped adults¹ (129 males, 113 females) were available for this study, representing the majority of the target Human Microbiome Project (HMP) cohort of 300 individuals. Adult subjects lacking evidence of disease were recruited based on a lengthy list of exclusion criteria; we will refer to them here as 'healthy', as defined by the consortium clinical sampling criteria (K. Aagaard *et al.*, manuscript submitted). Women were sampled at 18 body habitats, men at 15 (excluding three vaginal sites), distributed among five major body areas. Nine specimens were collected from the oral cavity and oropharynx: saliva; buccal mucosa (cheek), keratinized gingiva (gums), palate, tonsils, throat and tongue soft tissues, and supra- and subgingival dental plaque (tooth biofilm above and below the gum). Four skin specimens were collected from the two retroauricular creases (behind each ear) and the two antecubital fossae (inner elbows), and one specimen for the anterior nares (nostrils). A self-collected stool specimen represented the microbiota of the lower gastrointestinal tract, and three vaginal specimens were collected from the vaginal introitus, midpoint and posterior fornix. To evaluate within-subject stability of the microbiome, 131 individuals in these data were sampled at an additional time point (mean 219 days and s.d. 69 days after first sampling, range 35–404 days). After quality control, these specimens were used for 16S rRNA gene analysis via 454 pyrosequencing (abbreviated henceforth as 16S profiling, mean 5,408 and s.d. 4,605 filtered sequences per sample); to assess function, 681 samples were sequenced using paired-end Illumina shotgun metagenomic reads (mean 2.9 gigabases (Gb) and s.d. 2.1 Gb per sample)¹. More details on data generation are provided in related HMP publications¹ and in Supplementary Methods.

Microbial diversity of healthy humans

The diversity of microbes within a given body habitat can be defined as the number and abundance distribution of distinct types of organisms, which has been linked to several human diseases: low diversity in the gut to obesity and inflammatory bowel disease^{2,3}, for example, and high diversity in the vagina to bacterial vaginosis⁴. For this large study

involving microbiome samples collected from healthy volunteers at two distinct geographic locations in the United States, we have defined the microbial communities at each body habitat, encountering 81–99% of predicted genera and saturating the range of overall community configurations (Fig. 1, Supplementary Fig. 1 and Supplementary Table 1; see also Fig. 4). Oral and stool communities were especially diverse in terms of community membership, expanding prior observations⁵, and vaginal sites harboured particularly simple communities (Fig. 1a). This study established that these patterns of alpha diversity (within samples) differed markedly from comparisons between samples from the same habitat among subjects (beta diversity, Fig. 1b). For example, the saliva had among the highest median alpha diversities of operational taxonomic units (OTUs, roughly species level classification, see <http://hmpdacc.org/HMQCP>), but one of the lowest beta diversities—so although each individual's saliva was ecologically rich, members of the population shared similar organisms. Conversely, the antecubital fossae (skin) had the highest beta diversity but were intermediate in alpha diversity. The vagina had the lowest alpha diversity, with quite low beta diversity at the genus level but very high among OTUs due to the presence of distinct *Lactobacillus* spp. (Fig. 1b). The primary patterns of variation in community structure followed the major body habitat groups (oral, skin, gut and vaginal), defining as a result the complete range of population-wide between-subject variation in human microbiome habitats (Fig. 1c). Within-subject variation over time was consistently lower than between-subject variation, both in organismal composition and in metabolic function (Fig. 1d). The uniqueness of each individual's microbial community thus seems to be stable over time (relative to the population as a whole), which may be another feature of the human microbiome specifically associated with health.

No taxa were observed to be universally present among all body habitats and individuals at the sequencing depth employed here, unlike several pathways (Fig. 2 and Supplementary Fig. 2, see below), although several clades demonstrated broad prevalence and relatively abundant carriage patterns^{6,7}. Instead, as suggested by individually

*Lists of participants and their affiliations appear at the end of the paper.

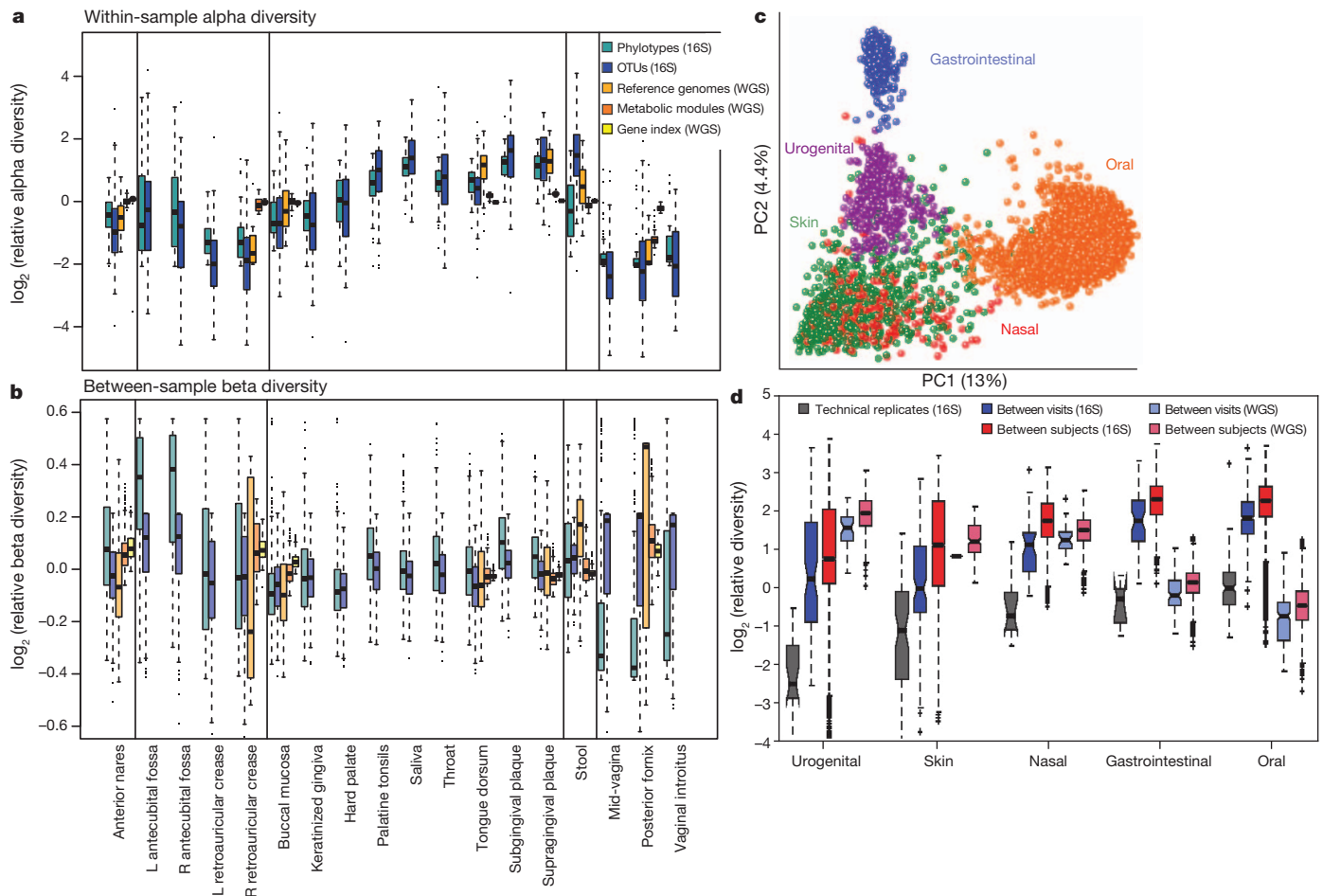


Figure 1 | Diversity of the human microbiome is concordant among measures, unique to each individual, and strongly determined by microbial habitat. **a**, Alpha diversity within subjects by body habitat, grouped by area, as measured using the relative inverse Simpson index of genus-level phylotypes (cyan), 16S rRNA gene OTUs (blue), shotgun metagenomic reads matched to reference genomes (orange), functional modules (dark orange), and enzyme families (yellow). The mouth generally shows high within-subject diversity and the vagina low diversity, with other habitats intermediate; variation among individuals often exceeds variation among body habitats. **b**, Bray–Curtis beta diversity among subjects by body habitat, colours as for **a**. Skin differs most between subjects, with oral habitats and vaginal genera more stable. Although

alpha- and beta-diversity are not directly comparable, changes in structure among communities (**a**) occupy a wider dynamic range than do changes within communities among individuals (**b**). **c**, Principal coordinates plot showing variation among samples demonstrates that primary clustering is by body area, with the oral, gastrointestinal, skin and urogenital habitats separate; the nares habitat bridges oral and skin habitats. **d**, Repeated samples from the same subject (blue) are more similar than microbiomes from different subjects (red). Technical replicates (grey) are in turn more similar; these patterns are consistent for all body habitats and for both phylogenetic and metabolic community composition. See previously described sample counts¹ for all comparisons.

focused studies^{2,3,5,8,9}, each body habitat in almost every subject was characterized by one or a few signature taxa making up the plurality of the community (Fig. 3). Signature clades at the genus level formed on average anywhere from 17% to 84% of their respective body habitats, completely absent in some communities (0% at this level of detection) and representing the entire population (100%) in others. Notably, less dominant taxa were also highly personalized, both among individuals and body habitats; in the oral cavity, for example, most habitats are dominated by *Streptococcus*, but these are followed in abundance by *Haemophilus* in the buccal mucosa, *Actinomyces* in the supragingival plaque, and *Prevotella* in the immediately adjacent (but low oxygen) subgingival plaque¹⁰.

Additional taxonomic detail of the human microbiome was provided by identifying unique marker sequences in metagenomic data¹¹ (Fig. 3a) to complement 16S profiling (Fig. 3b). These two profiles were typically in close agreement (Supplementary Fig. 3), with the former in some cases offering more specific information on members of signature genera differentially present within habitats (for example, vaginal *Prevotella amnii* and gut *Prevotella copri*) or among individuals (for example, vaginal *Lactobacillus* spp.). One application of this specificity was to confirm the absence of NIAID (National Institute of

Allergy and Infectious Diseases) class A–C pathogens above 0.1% abundance (aside from *Staphylococcus aureus* and *Escherichia coli*) from the healthy microbiome, but the near-ubiquity and broad distribution of opportunistic ‘pathogens’ as defined by PATRIC¹². Canonical pathogens including *Vibrio cholerae*, *Mycobacterium avium*, *Campylobacter jejuni* and *Salmonella enterica* were not detected at this level of sensitivity. *Helicobacter pylori* was found in only two stool samples, both at <0.01%, and *E. coli* was present at >0.1% abundance in 15% of stool microbiomes (>0% abundance in 61%). Similar species-level observations were obtained for a small subset of stool samples with 454 pyrosequencing metagenomics data using PhylOTU^{13,14}. In total 56 of 327 PATRIC pathogens were detected in the healthy microbiome (at >1% prevalence of >0.1% abundance, Supplementary Table 2), all opportunistic and, strikingly, typically prevalent both among hosts and habitats. The latter is in contrast to many of the most abundant signature taxa, which were usually more habitat-specific and variable among hosts (Fig. 3a, b). This overall absence of particularly detrimental microbes supports the hypothesis that even given this cohort’s high diversity, the microbiota tend to occupy a range of configurations in health distinct from many of the disease perturbations studied to date^{3,15}.

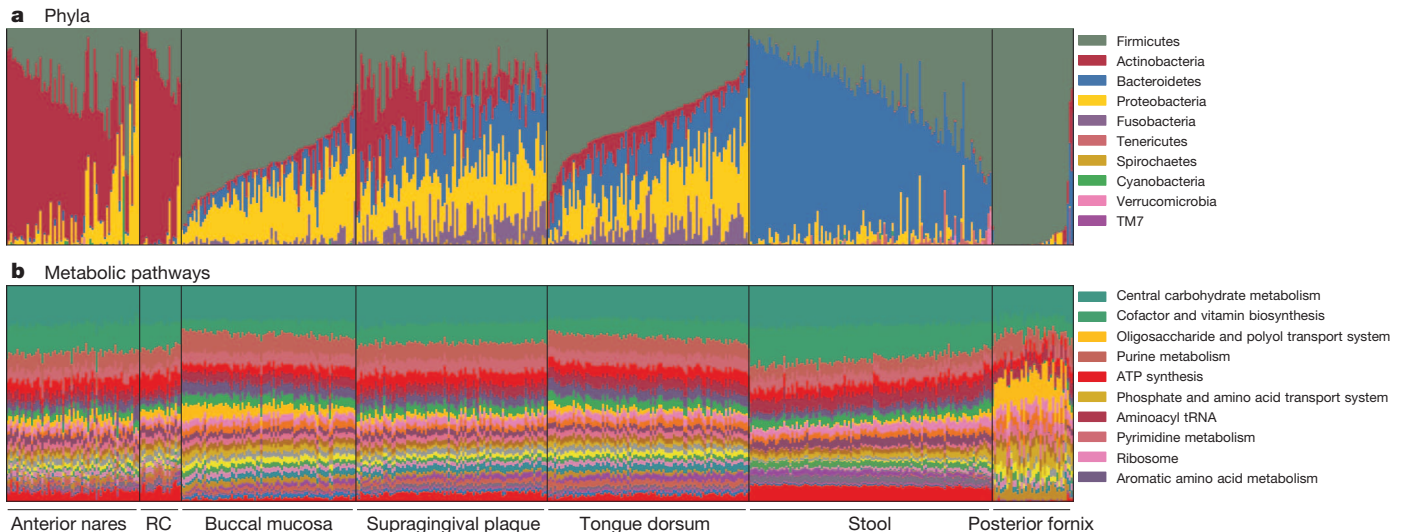


Figure 2 | Carriage of microbial taxa varies while metabolic pathways remain stable within a healthy population. a, b, Vertical bars represent microbiome samples by body habitat in the seven locations with both shotgun and 16S data; bars indicate relative abundances colored by microbial phyla from binned OTUs (a) and metabolic modules (b). Legend indicates most abundant phyla/pathways by average within one or more body habitats; RC,

retroauricular crease. A plurality of most communities' memberships consists of a single dominant phylum (and often genus; see Supplementary Fig. 2), but this is universal neither to all body habitats nor to all individuals. Conversely, most metabolic pathways are evenly distributed and prevalent across both individuals and body habitats.

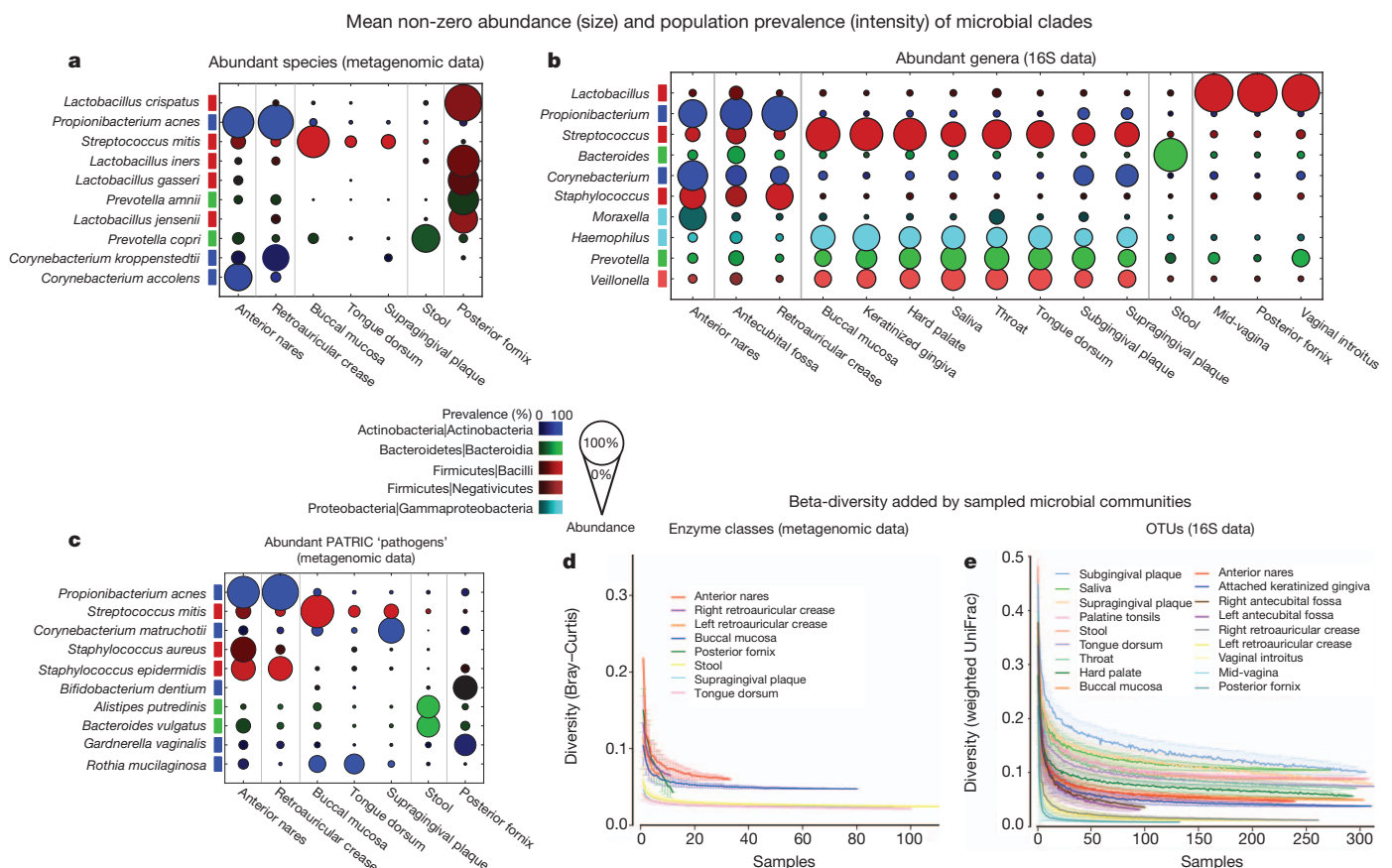


Figure 3 | Abundant taxa in the human microbiome that have been metagenomically and taxonomically well defined in the HMP population. a–c, Prevalence (intensity, colour denoting phylum/class) and abundance when present (size) of clades in the healthy microbiome. The most abundant metagenomically-identified species (a), 16S-identified genera (b) and PATRIC¹² pathogens (metagenomic) (c) are shown. d, e, The population size

and sequencing depths of the HMP have well defined the microbiome at all assayed body sites, as assessed by saturation of added community metabolic configurations (rarefaction of minimum Bray–Curtis beta-diversity of metagenomic enzyme class abundances to nearest neighbour, inter-quartile range over 100 samples) (d) and phylogenetic configurations (minimum 16S OTU weighted UniFrac distance to nearest neighbour) (e).

Carriage of specific microbes

Inter-individual variation in the microbiome proved to be specific, functionally relevant and personalized. One example of this is illustrated by the *Streptococcus* spp. of the oral cavity. The genus dominates the oropharynx¹⁶, with different species abundant within each sampled body habitat (see <http://hmpdacc.org/HMSMCP>) and, even at the species level, marked differences in carriage within each habitat among individuals (Fig. 4a). As the ratio of pan- to core-genomes is high in many human-associated microbes¹⁷, this variation in abundance could be due to selective pressures acting on pathways differentially present among *Streptococcus* species or strains (Fig. 4b). Indeed, we observed extensive strain-level genomic variation within microbial species in this population, enriched for host-specific structural variants around genomic islands (Fig. 4c). Even with respect to the single *Streptococcus mitis* strain B6, gene losses associated with these events were common,

for example differentially eliminating *S. mitis* carriage of the V-type ATPase or choline binding proteins cbp6 and cbp12 among subsets of the host population (Fig. 4d). These losses were easily observable by comparison to reference isolate genomes, and these initial findings indicate that microbial strain- and host-specific gene gains and polymorphisms may be similarly ubiquitous.

Other examples of functionally relevant inter-individual variation at the species and strain levels occurred throughout the microbiome. In the gut, *Bacteroides fragilis* has been shown to prime T-cell responses in animal models via the capsular polysaccharide A¹⁸, and in the HMP stool samples this taxon was carried at a level of at least 0.1% in 16% of samples (over 1% abundance in 3%). *Bacteroides thetaiotaomicron* has been studied for its effect on host gastrointestinal metabolism¹⁹ and was likewise common at 46% prevalence. On the skin, *S. aureus*, of particular interest as the cause of methicillin-resistant

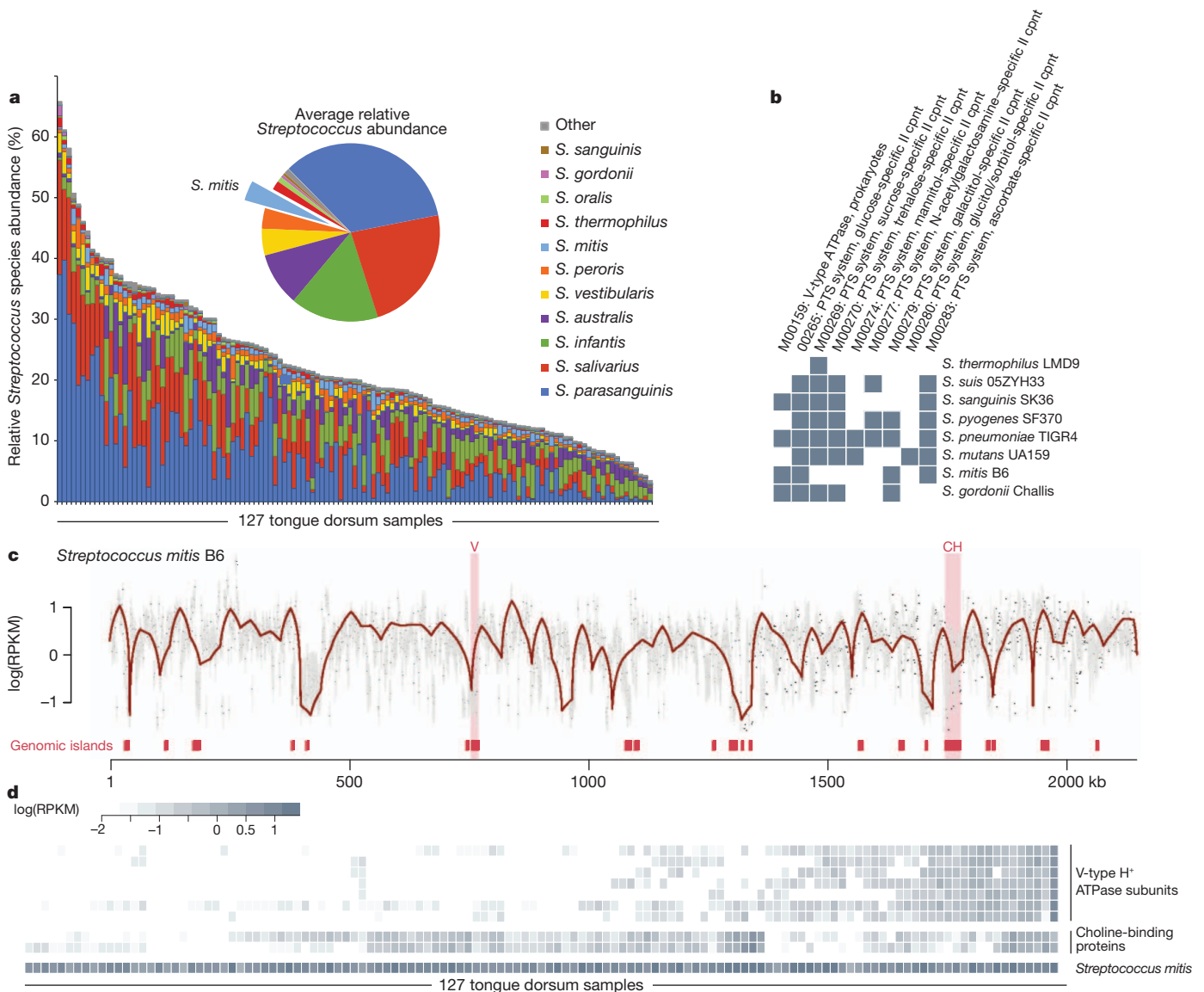


Figure 4 | Microbial carriage varies between subjects down to the species and strain level. Metagenomic reads from 127 tongue samples spanning 90 subjects were processed with MetaPhlAn to determine relative abundances for each species. **a**, Relative abundances of 11 distinct *Streptococcus* spp. In addition to variation in broader clades (see Fig. 2), individual species within a single habitat demonstrate a wide range of compositional variation. Inset illustrates average tongue sample composition. **b**, Metabolic modules present/absent (grey/white) in KEGG²⁴ reference genomes of tongue streptococci denote selected areas of strain-specific functional differentiation. cpnt, component.

c, Comparative genomic coverage for the single *Streptococcus mitis* B6 strain. Grey dots are median reads per kilobase per million reads (RPKM) for 1-kb windows, grey bars are the 25th to 75th percentiles across all samples, red line the LOWESS-smoothed average. Red bars at the bottom highlight predicted genomic islands²⁷. Large, discrete, and highly variable islands are commonly under-represented. **d**, Two islands are highlighted, V (V-type H⁺ ATPase subunits I, K, E, C, F, A and B) and CH (choline-binding proteins cbp6 and cbp12), indicating functional cohesion of strain-specific gene loss within individual human hosts.

S. aureus (MRSA) infections, had 29% nasal and 4% skin carriage rates, roughly as expected²⁰. Close phylogenetic relatives such as *Staphylococcus epidermidis* (itself considered commensal) were, in contrast, universal on the skin and present in 93% of nares samples, and at the opposite extreme *Pseudomonas aeruginosa* (a representative Gram-negative skin pathogen) was completely absent from both body habitats (0% at this level of detection). These and the data above suggest that the carriage pattern of some species in the human microbiome may be analogous to genetic traits, where recessive alleles of modest risk are maintained in a population. In the case of the human microbiome, high-risk pathogens remain absent, whereas species that pose a modest degree of risk also seem to be stably maintained in this ecological niche.

Finally, microorganisms within and among body habitats exhibited relationships suggestive of driving physical factors such as oxygen, moisture and pH, host immunological factors, and microbial interactions such as mutualism or competition²¹ (Supplementary Fig. 4). Both overall community similarity and microbial co-occurrence and co-exclusion across the human microbiome grouped the 18 body habitats together into four clusters corresponding to the five target body areas (Supplementary Fig. 4a, b). There was little distinction among different vaginal sites, with *Lactobacillus* spp. dominating all three and correlating in abundance. However, *Lactobacillus* varied inversely with the Actinobacteria and Bacteroidetes (see Supplementary Fig. 4c and Figs 2 and 3), as also observed in a previous cohort⁹. Gut microbiota relationships primarily comprised inverse associations with the *Bacteroides*, which ranged from dominant in some subjects to a minority in others who carried a greater diversity of Firmicutes. A similar progression was evident in the skin communities, dominated by one of *Staphylococcus* (phylum Firmicutes), *Propionibacterium*, or *Corynebacterium* (both phylum Actinobacteria), with a continuum of oral organisms (for example, *Streptococcus*) appearing in nares communities (Supplementary Fig. 4c). These observations suggest that microbial community structure in these individuals may sometimes occupy discrete configurations and under other circumstances vary continuously, a topic addressed in more detail by several HMP investigations (ref. 6 and unpublished results). An individual's location within such configurations is indicative of current microbial carriage (including pathogens) and of the community's ability to resist future pathogen acquisition or dysbiosis; it may thus prove to be associated with disease susceptibility or other phenotypic characteristics.

Microbiome metabolism and function

As the first study to include both marker gene and metagenomic data across body habitats from a large human population, we additionally assessed the ecology of microbial metabolic and functional pathways in these communities. We reconstructed the relative abundances of pathways in community metagenomes²², which were much more constant and evenly diverse than were organismal abundances (Fig. 2b, see also Fig. 1), confirming this as an ecological property of the entire human microbiome². We were likewise able to determine for the first time that taxonomic and functional alpha diversity across microbial communities significantly correlate (Spearman of inverse Simpson's $r = 0.60$, $P = 3.6 \times 10^{-67}$, $n = 661$), the latter within a more proscribed range of community configurations (Supplementary Fig. 5).

Unlike microbial taxa, several pathways were ubiquitous among individuals and body habitats. The most abundant of these 'core' pathways include the ribosome and translational machinery, nucleotide charging and ATP synthesis, and glycolysis, and reflect the basics of host-associated microbial life. Also in contrast to taxa, few pathways were highly variable among subjects within any body habitat; exceptions included the Sec (orally, pathway relative abundance s.d. = 0.0052; total mean of oral standard deviations = 0.0011 with s.d. = 0.0016) and Tat (globally, pathway s.d. = 0.0055; mean of

global standard deviations = 0.0023 with s.d. = 0.0033) secretion systems, indicating a high degree of host-microbe and microbe-microbe interactions in the healthy human microbiota. This high variability was particularly present in the oral cavity; for phosphate, mono- and di-saccharide, and amino acid transport in the mucosa; and also for lipopolysaccharide biosynthesis and spermidine/putrescine synthesis and transport on the plaque and tongue (<http://hmpdacc.org/HMMRC>). The stability and high metagenomic abundance of this housekeeping 'core' contrasts with the greater variability and lower abundance of niche-specific functionality in rare but consistently present pathways; for example, spermidine biosynthesis, methionine degradation and hydrogen sulphide production, all examples highly prevalent in gastrointestinal body sites (non-zero in >92% of samples) but at very low abundance (median relative abundance < 0.0052). This 'long tail' of low-abundance genes and pathways also probably encodes much of the uncharacterized biomolecular function and metabolism of these metagenomes, the expression levels of which remain to be explored in future metatranscriptomic studies.

Protein families showed diversity and prevalence trends similar to those of full pathways, ranging from maxima of only ~16,000 unique families per community in the vagina to almost 400,000 in the oral cavity (Fig. 1a, b; <http://hmpdacc.org/HMGI>). A remarkable fraction of these families were indeed functionally uncharacterized, including those detected by read mapping, with a minimum in the oral cavity (mean 58% s.d. 6.8%) and maximum in the nares (mean 77% s.d. 11%). Likewise, many genes annotated from assemblies could not be assigned a metabolic function, with a minimum in the vagina (mean 78% s.d. 3.4%) and maximum in the gut (mean 86% s.d. 0.9%). The latter range did not differ substantially by body habitat and is in close agreement with previous comprehensive gene catalogues of the gut metagenome³. Taken together with the microbial variation observed above throughout the human microbiome, functional variation among individuals might indicate pathways of particular importance in maintaining community structure in the face of personalized immune, environmental or dietary exposures among these subjects. Determining the functions of uncharacterized core and variable protein families will be especially essential in understanding role of the microbiota in health and disease.

Correlations with host phenotype

We finally examined relationships associating both clades and metabolism in the microbiota with host properties such as age, gender, body mass index (BMI), and other available clinical metadata (Fig. 5 and Supplementary Table 3). Using a sparse multivariate model, 960 microbial, enzymatic or pathway abundances were significantly associated with one or more of 15 subject phenotype and sample metadata features. A wide variety of taxa, gene families and metabolic pathways were differentially distributed with subject ethnicity at every body habitat (Fig. 5a), representing the phenotype with the greatest number (266 at false discovery rate (FDR) $q < 0.2$) of total associations with the microbiome. Vaginal pH has also been observed to correlate with microbiome composition⁹, and we detected in this population both the expected reduction in *Lactobacillus* at high pH and a corresponding increase in metabolic diversity (Fig. 5b). Intriguingly, and not previously observed, subject age was most associated with a collection of highly differential metagenomically encoded pathways on the skin (Fig. 5c), as well as shifts in skin clades including retroauricular Firmicutes ($P = 1.0 \times 10^{-4}$, $q = 0.033$). The examples of associations with ethnicity and vaginal pH are among the strongest associations with the microbiome, however, and most correlates (for example, with subject BMI, Fig. 5d) are more representatively modest. This lower degree of correlation held for most available biometrics (gender, temperature, blood pressure, etc.), with even the most significant associations possessing generally low effect sizes and considerable unexplained variance. We conclude that most variation in the human microbiome is not well explained by these phenotypic

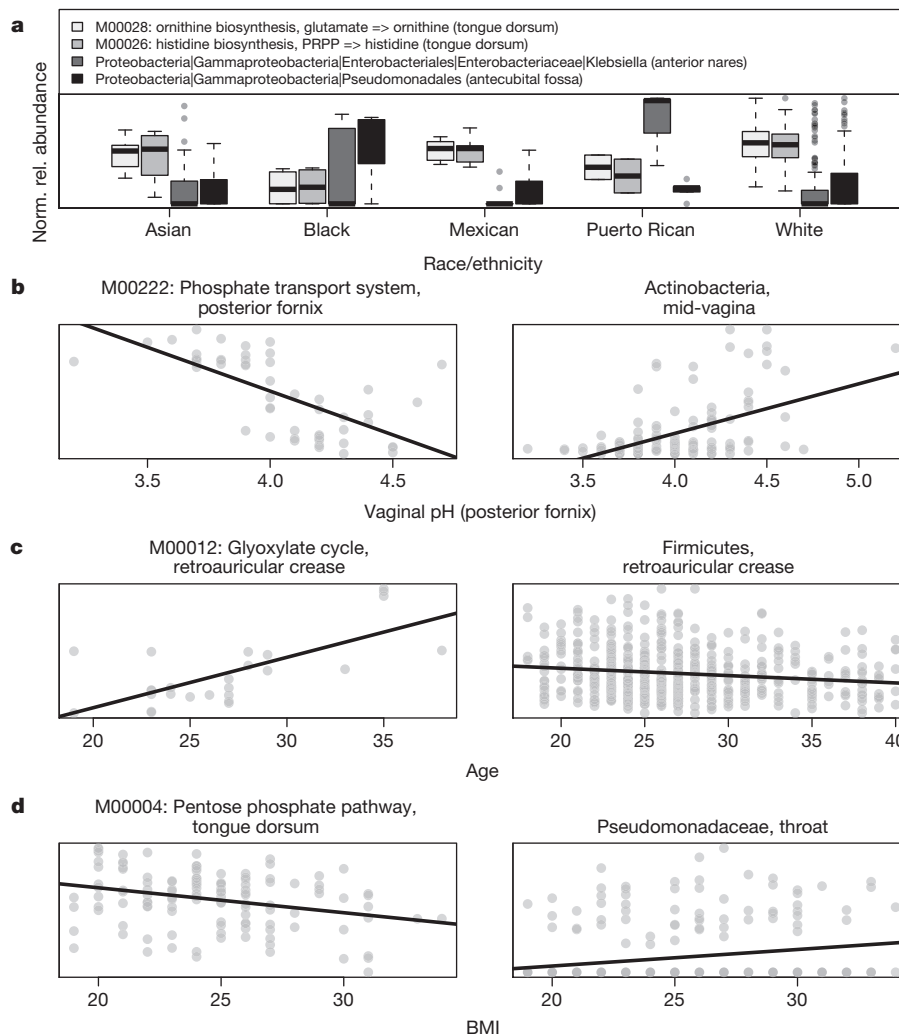


Figure 5 | Microbial community membership and function correlates with host phenotype and sample metadata. **a–d**, The pathway and clade abundances most significantly associated (all FDR $q < 0.2$) using a multivariate linear model with subject race or ethnicity (**a**), vaginal posterior fornix pH (**b**), subject age (**c**) and BMI (**d**). Scatter plots of samples are shown with lines

indicating best simple linear fit. Race/ethnicity and vaginal pH are particularly strong associations; age and BMI are more representative of typically modest phenotypic associations (Supplementary Table 3), suggesting that variation in the healthy microbiota may correspond to other host or environmental factors.

metadata, and other potentially important factors such as short- and long-term diet, daily cycles, founder effects such as mode of delivery, and host genetics should be considered in future analyses.

Conclusions

This extensive sampling of the human microbiome across many subjects and body habitats provides an initial characterization of the normal microbiota of healthy adults in a Western population. The large sample size and consistent sampling of many sites from the same individuals allows for the first time an understanding of the relationships among microbes, and between the microbiome and clinical parameters, that underpin the basis for individual variation—variation that may ultimately be critical for understanding microbiome-based disorders. Clinical studies of the microbiome will be able to leverage the resulting extensive catalogues of taxa, pathways and genes¹, although they must also still include carefully matched internal controls. The uniqueness of each individual's microbiome even in this reference population argues for future studies to consider prospective within-subjects designs where possible. The HMP's unique combination of organismal and functional data across body habitats, encompassing both 16S and metagenomic profiling, together with detailed characterization of each subject, has allowed us and subsequent studies to move beyond the observation of

variability in the human microbiome to ask how and why these microbial communities vary so extensively.

Many details remain for further work to fill in, building on this reference study. How do early colonization and lifelong change vary among body habitats? Do epidemiological patterns of transmission of beneficial or harmless microbes mirror patterns of transmission of pathogens? Which co-occurrences among microbes reflect shared response to the environment, as opposed to competitive or mutualistic interactions? How large a role does host immunity or genetics play in shaping patterns of diversity, and how do the patterns observed in this North American population compare to those around the world? Future studies building on the gene and organism catalogues established by the Human Microbiome Project, including increasingly detailed investigations of metatranscriptomes and metaproteomes, will help to unravel these open questions and allow us to more fully understand the links between the human microbiome, health and disease.

METHODS SUMMARY

Microbiome samples were collected from up to 18 body sites at one or two time points from 242 individuals clinically screened for absence of disease (K. Aagaard *et al.*, manuscript submitted). Samples were subjected to 16S ribosomal RNA gene pyrosequencing (454 Life Sciences), and a subset were shotgun-sequenced for metagenomics using the Illumina GAIIx platform¹. 16S data processing and

diversity estimates were performed using QIIME²³, and metagenomic data were taxonomically profiled using MetaPhlAn¹¹, metabolically profiled by HUMAnN²², and assembled for gene annotation and clustering into a unique catalogue¹. Potential pathogens were identified using the PATRIC database¹², isolate reference genome annotations drawn from KEGG²⁴, and reference genome mapping performed by BWA²⁵ to a reduced set of genomes to which short reads could be matched²⁶. Microbial associations were assessed by similarity measures accounting for compositionality²⁷, and phenotypic association testing was performed in R. All data and additional protocol details are available at <http://hmpdacc.org>. Full methods accompany this paper in the Supplementary Information.

Received 2 November 2011; accepted 16 May 2012.

1. The Human Microbiome Project Consortium. A framework for human microbiome research. *Nature* <http://dx.doi.org/10.1038/nature11209> (this issue).
2. Turnbaugh, P. J. *et al.* A core gut microbiome in obese and lean twins. *Nature* **457**, 480–484 (2009).
3. Qin, J. *et al.* A human gut microbial gene catalogue established by metagenomic sequencing. *Nature* **464**, 59–65 (2010).
4. Fredricks, D. N., Fiedler, T. L. & Marrazzo, J. M. Molecular identification of bacteria associated with bacterial vaginosis. *N. Engl. J. Med.* **353**, 1899–1911 (2005).
5. Costello, E. K. *et al.* Bacterial community variation in human body habitats across space and time. *Science* **326**, 1694–1697 (2009).
6. Huse, S., Ye, Y., Zhou, Y. & Fodor, A. A core human microbiome as viewed through 16S rRNA sequences clusters. *PLoS ONE* <http://dx.doi.org/10.1371/journal.pone.0034242> (14 June 2012).
7. Li, K., Bihan, M., Yooseph, S. & Methe, B. A. Analyses of the microbial diversity across the human microbiome. *PLoS ONE* <http://dx.doi.org/10.1371/journal.pone.0032118> (14 June 2012).
8. Grice, E. A. *et al.* Topographical and temporal diversity of the human skin microbiome. *Science* **324**, 1190–1192 (2009).
9. Ravel, J. *et al.* Vaginal microbiome of reproductive-age women. *Proc. Natl Acad. Sci. USA* **108** (Suppl 1), 4680–4687 (2011).
10. Segata, N. *et al.* Composition of the adult digestive tract microbiome based on seven mouth surfaces, tonsils, throat and stool samples. *Genome Biol.* **13**, R42 (2012).
11. Segata, N. *et al.* Efficient metagenomic microbial community profiling using unique clade-specific marker genes. *Nature Methods* <http://dx.doi.org/10.1038/nmeth.2066> (2012).
12. Gillespie, J. J. *et al.* PATRIC: the comprehensive bacterial bioinformatics resource with a focus on human pathogenic species. *Infect. Immun.* **79**, 4286–4298 (2011).
13. Sharpton, T. J. *et al.* PhylOTU: a high-throughput procedure quantifies microbial community diversity and resolves novel taxa from metagenomic data. *PLoS Comput. Biol.* **7**, e1001061 (2011).
14. Wylie, K. M. *et al.* Novel bacterial taxa in the human microbiome. *PLoS ONE* <http://dx.doi.org/10.1371/journal.pone.0035229> (14 June 2012).
15. Sokol, H. *et al.* *Faecalibacterium prausnitzii* is an anti-inflammatory commensal bacterium identified by gut microbiota analysis of Crohn disease patients. *Proc. Natl Acad. Sci. USA* **105**, 16731–16736 (2008).
16. Aas, J. A., Paster, B. J., Stokes, L. N., Olsen, I. & Dewhirst, F. E. Defining the normal bacterial flora of the oral cavity. *J. Clin. Microbiol.* **43**, 5721–5732 (2005).
17. Medini, D. *et al.* Microbiology in the post-genomic era. *Nature Rev. Microbiol.* **6**, 419–430 (2008).
18. Mazmanian, S. K., Round, J. L. & Kasper, D. L. A microbial symbiosis factor prevents intestinal inflammatory disease. *Nature* **453**, 620–625 (2008).
19. Goodman, A. L. *et al.* Identifying genetic determinants needed to establish a human gut symbiont in its habitat. *Cell Host Microbe* **6**, 279–289 (2009).
20. Kuehnert, M. J. *et al.* Prevalence of *Staphylococcus aureus* nasal colonization in the United States, 2001–2002. *J. Infect. Dis.* **193**, 172–179 (2006).
21. Faust, K. *et al.* Microbial co-occurrence relationships in the human microbiome. *PLoS Comput. Biol.* (in the press).
22. Abubucker, S. *et al.* Metabolic reconstruction for metagenomic data and its application to the human microbiome. *PLoS Comput. Biol.* <http://dx.doi.org/10.1371/journal.pcbi.1002358> (14 June 2012).
23. Caporaso, J. G. *et al.* QIIME allows analysis of high-throughput community sequencing data. *Nature Methods* **7**, 335–336 (2010).
24. Kanehisa, M., Goto, S., Furumichi, M., Tanabe, M. & Hirakawa, M. KEGG for representation and analysis of molecular networks involving diseases and drugs. *Nucleic Acids Res.* **38**, D355–D360 (2010).
25. Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows–Wheeler transform. *Bioinformatics* **26**, 589–595 (2010).
26. Giannoukos, G. *et al.* Efficient and robust RNA-seq process for cultured bacteria and complex community transcriptomes. *Genome Biol.* **13**, R23 (2012).
27. Langille, M. G. & Brinkman, F. S. IslandViewer: an integrated interface for computational identification and visualization of genomic islands. *Bioinformatics* **25**, 664–665 (2009).

Supplementary Information is linked to the online version of the paper at www.nature.com/nature.

Acknowledgements The Consortium would like to thank our external scientific advisory board: R. Blumberg, J. Davies, R. Holt, P. Ossorio, F. Ouellette, G. Schoolnik and A. Williamson. We would also like to thank our collaborators throughout the

International Human Microbiome Consortium, particularly the investigators of the MetaHIT project, for advancing human microbiome research. Data repository management was provided by the National Center for Biotechnology Information and the Intramural Research Program of the NIH National Library of Medicine. We appreciate the participation of the individuals from the Saint Louis, Missouri, and Houston, Texas areas who made this study possible. This research was supported in part by National Institutes of Health grants U54HG004969 to B.W.B.; U54HG003273 to R.A.G.; U54HG004973 to R.A.G., S.K.H. and J.F.P.; U54HG003067 to E.S.Lander; U54AI084844 to K.E.N.; N01AI30071 to R.L. Strausberg; U54HG004968 to G.M.W.; U01HG004866 to O.R.W.; U54HG003079 to R.K.W.; R01HG005969 to C.H.; R01HG004872 to R.K.; R01HG004885 to M.P.; R01HG005975 to P.D.S.; R01HG004908 to Y.Y.; R01HG004900 to M.K.Cho and P. Sankar; R01HG005171 to D.E.H.; R01HG004853 to A.L.M.; R01HG004856 to R.R.; R01HG004877 to R.R.S. and R.F.; R01HG005172 to P. Spicer; R01HG004857 to M.P.; R01HG004906 to T.M.S.; R21HG005811 to E.A.V.; M.J.B. was supported by UH2AR057506; G.A.B. was supported by UH2AI083263 and UH3AI083263 (G.A.B., C. N. Cornelissen, L. K. Eaves and J. F. Strauss); S.M.H. was supported by UH3DK083993 (V. B. Young, E. B. Chang, F. Meyer, T. M. S., M. L. Sogin, J. M. Tiedje); K.P.R. was supported by UH2DK083990 (J. V.); J.A.S. and H.H.K. were supported by UH2AR057504 and UH3AR057504 (J.A.S.); DP2OD001500 to K.M.A.; N01HG62088 to the Coriell Institute for Medical Research; U01DE016937 to F.E.D.; S.K.H. was supported by R01DE0202098 and R01DE021574 (S.K.H. and H. Li); J.I. was supported by R21CA139193 (J.I. and D. S. Michaud); K.P.L. was supported by P30DE020751 (D. J. Smith); Army Research Office grant W911NF-11-1-0473 to C.H.; National Science Foundation grants NSF DBI-1053486 to C.H. and NSF IIS-0812111 to M.P.; The Office of Science of the US Department of Energy under Contract No. DE-AC02-05CH11231 for P.S. C.; LANL Laboratory-Directed Research and Development grant 20100034DR and the US Defense Threat Reduction Agency grants B104153I and B084531I to P.S.C.; Research Foundation – Flanders (FWO) grant to K.F. and J.Raes; R.K. is an HHMI Early Career Scientist; Gordon & Betty Moore Foundation funding and institutional funding from the J. David Gladstone Institutes to K.S.P.; A.M.S. was supported by fellowships provided by the Rackham Graduate School and the NIH Molecular Mechanisms in Microbial Pathogenesis Training Grant T32AI007528; a Crohn's and Colitis Foundation of Canada Grant in Aid of Research to E.A.V.; 2010 IBM Faculty Award to K.C.W.; analysis of the HMP data was performed using National Energy Research Scientific Computing resources, the BluBioU Computational Resource at Rice University.

Author Contributions Principal investigators: B.W.B., R.A.G., S.K.H., B.A.M., K.E.N., J.F.P., G.M.W., O.W., R.K.W. Manuscript preparation: D.G., C.H., R.K., O.W. Funding agency management: C.C.B., T.B., V.R.B., J.L.C., S.C., C.D., V.D.F., C.G., M.Y.G., R.D.L., J.M., P.M., J.P., L.M.P., J.A.S., L.W., C.W., K.A.W. Project leadership: S.A., J.H.B., B.W.B., A.T.C., H.H.C., A.M.E., M.G.F., R.S.F., D.G., M.G.G., K.H., S.K.H., C.H., E.A.L., R.M., V.M., J.C.M., B.A.M., M.M., D.M.M., K.E.N., J.F.P., E.J.S., J.V., G.M.W., O.W., A.M.W., K.C.W., J.R.W., S.K.Y., Q.Z. Analysis preparation for manuscript: J.C.C., K.F., D.G., A.G., K.H.H., C.H., R.K., D.K., H.H.K., O.K., K.P.L., R.E.L., J.R., J.F.S., P.D.S., N.S. Data release: L.A., T.B., I.A.C., K.C., H.H.C., N.J.D., D.J.D., A.M.E., V.M.F., L.F., J.M.G., S.G., S.K.H., M.E.H., C.J., V.J., C.K., A.A.M., V.M.M., T.M., M.M., D.M.M., J.O., K.P., J.F.P., C.P., X.Q., R.K.S., N.S., I.S., E.J.S., D.V.W., O.W., K.W., K.C.W., C.Y., B.P.Y., Q.Z. Methods and research development: S.A., H.M.A., M.B., D.M.C., A.M.E., R.L.E., M.F., S.F., M.G.F., D.C.F., D.G., G.G., B.J.H., S.K.H., M.E.H., W.A.K., N.L., K.L., V.M., E.R.M., B.A.M., M.M., D.M.M., C.N., J.F.P., M.E.P., X.Q., M.C.R., C.R., E.J.S., S.M.S., D.G.T., D.V.W., G.M.W., Y.W., K.A.W., S.Y., B.P.Y., S.K.Y., Q.Z. DNA sequence production: S.A., E.A., T.A., T.B., C.J.B., D.A.B., K.D.D., S.P.D., A.M.E., R.L.E., C.N.F., S.F., C.C.F., L.L.F., R.S.F., B.H., S.K.H., M.E.H., V.J., C.L.K., S.L.L., N.L., L.L., D.M.M., I.N., C.N., M.O., J.F.P., X.Q., J.G.R., Y.R., M.C.R., D.V.W., Y.W., B.P.Y., Y.Z. Clinical sample collection: K.M.A., M.A.C., W.M.D., L.L.F., N.G., H.A.H., E.L.H., J.A.K., W.A.K., T.M., A.L.M., P.M., S.M.P., J.F.P., G.A.S., J.V., M.A.W., G.M.W. Body site experts: K.M.A., E.A.V., G.A., L.B., M.J.B., C.C.D., F.E.D., L.F., J.I., J.A.K., S.K.H., H.H.K., K.P.L., P.J.M., J. Ravel, T.M.S., J.A.S., J.D.S., J.V. Ethical, legal and social implications: R.M.F., D.E.H., W.A.K., N.B.K., C.M.L., A.L.M., R.R., P. Sankar, R.R.S., P. Spicer, L.Z. Strain management: E.A.V., J.H.B., I.A.C., K.C., S.W.C., H.H.C., T.Z.D., A.S.D., A.M.E., M.G.F., M.G.G., S.K.H., V.J., N.C.K., S.L.L., L.L., K.L., E.A.L., V.M.M., B.A.M., D.M.M., K.E.N., I.N., I.P., L.S., E.J.S., C.M.T., M.T., D.V.W., G.M.W., A.M.W., Y.W., K.M.W., B.P.Y., L.Z., Y.Z. 16S data analysis: K.M.A., E.J.A., G.L.A., C.A.A., M.B., B.W.B., J.P.B., G.A.B., S.R.C., S.C., J.C., T.Z.D., F.E.D., E.D., A.M.E., R.C.E., K.F., M.F., A.A.F., J.F., H.G., D.G., B.J.H., T.A.H., S.M.H., C.H., J.I., J.K.J., S.T.K., S.K.H., R.K., H.H.K., O.K., P.S.L., R.E.L., K.L., C.A.L., D.M., B.A.M., K.A.M., M.M., M.P., J.F.P., M.P., K.S.P., X.Q., J. Raes, K.P.R., M.C.R., B.R., J.F.S., P.D.S., T.M.S., N.S., J.A.S., W.D.S., T.J.S., C.S.S., E.J.S., R.M.T., J.V., T.A.V., Z.W., D.V.W., G.M.W., J.R.W., K.M.W., Y.Y., S.Y., Y.Z. Shotgun data processing and alignments: C.J.B., J.C.C., E.D., D.G., A.G., M.E.H., H.J., D.K., K.C.K., C.L.K., Y.L., J.C.M., B.A.M., M.M., D.M.M., J.O., J.F.P., X.Q., J.G.R., R.K.S., N.U.S., I.S., E.J.S., G.G.S., S.M.S., J.W., Z.W., G.M.W., O.W., K.C.W., T.W., S.K.Y., L.Z. Assembly: H.M.A., C.J.B., P.S.C., L.C., Y.D., S.P.D., M.G.F., M.E.H., H.J., S.K., B.L., Y.L., C.L., J.C.M., J.M.M., J.R.M., P.J.M., M.M., J.F.P., M.P., M.E.P., X.Q., M.R., R.K.S., M.S., D.D.S., G.G.S., S.M.S., C.M.T., T.J.T., W.W., G.M.W., K.C.W., L.Y., Y.Y., S.K.Y., L.Z. Annotation: O.O.A., V.B., C.J.B., I.A.C., A.T.C., K.C., H.H.C., A.S.D., M.G.G., J.M.G., J.G., A.G., S.G., B.J.H., K.H., S.K.H., C.H., H.J., N.C.K., R.M., V.M.M., K.M., T.M., M.M., J.O., K.P., M.P., X.Q., N.S., E.J.S., G.G.S., S.M.S., M.T., G.M.W., K.C.W., J.R.W., C.Y., S.K.Y., Q.Z., L.Z., W.G.S. Metabolic reconstruction: S.A., B.L.C., J.G., C.H., J.I., B.A.M., M.M., B.R., A.M.S., N.S., M.T., G.M.W., S.Y., Q.Z., J.D.Z.

Author Information All data used in this study is available from the Human Microbiome Project Data Analysis and Coordination Center at <http://hmpdacc.org> and from the NCBI. Reprints and permissions information is available at www.nature.com/reprints. This paper is distributed under the terms of the Creative Commons Attribution-Non-Commercial-Share Alike licence, and is freely available to all readers at www.nature.com/nature. The authors declare no competing financial interests. Readers are welcome to comment on the online version of this article at www.nature.com/nature. Correspondence and requests for materials should be addressed to C.H. (chuttenh@hsph.harvard.edu).

Curtis Huttenhower^{1,2*}, Dirk Gevers^{2*}, Rob Knight^{3,4}, Sahar Abubucker⁵, Jonathan H. Badger⁶, Asif T. Chinwalla⁵, Heather H. Creasy⁷, Ashlee M. Earl², Michael G. FitzGerald², Robert S. Fulton⁵, Michelle G. Giglio⁷, Kymberlie Hallsworth-Pepin⁵, Elizabeth A. Lobos⁵, Ramana Madupu⁶, Vincent Martin⁵, John C. Martin⁵, Makedonka Mitreva⁵, Donna M. Muzny⁸, Erica J. Sodergren⁵, James Versalovic^{9,10}, Aye M. Wollam⁵, Kim C. Worley⁸, Jennifer R. Wortman⁷, Sarah K. Young², Qiangdong Zeng², Kjersti M. Aagaard¹¹, Olukemi O. Abolude⁷, Emma Allen-Vercos¹², Eric J. Alm^{13,2}, Lucia Alvarado², Gary L. Andersen¹⁴, Scott Anderson², Elizabeth Appelbaum⁵, Harindra M. Arachchi², Gary Armitage¹⁵, Cesar A. Arze⁷, Tulin Ayvaz¹⁶, Carl C. Baker¹⁷, Lisa Begg¹⁸, Tsegahiwot Belachew¹⁹, Veena Bhongardien², Monika Bihan⁶, Martin J. Blaser²⁰, Toby Bloom², Vivien Bonazzi²¹, J. Paul Brooks^{22,23}, Gregory A. Buck^{23,24}, Christian J. Buhay⁸, Dana A. Busam⁶, Joseph L. Campbell^{21,19}, Shane R. Canon²⁵, Brandi L. Cantarel⁷, Patrick S. G. Chain^{26,27}, I-Min A. Chen²⁸, Lei Chen⁵, Shaila Chhibba²¹, Ken Chu²⁸, Dawn M. Ciulla², Jose C. Clemente³, Sandra W. Clifton⁵, Sean Conlan⁷⁹, Jonathan Crabtree⁷, Mary A. Cutting²⁹, Noam J. Davidovics⁷, Catherine C. Davis³⁰, Todd Z. DeSantis³¹, Carolyn Deal¹⁹, Kimberley D. Delehaunty³, Floyd E. Dewhirst^{32,33}, Elena Deych³⁴, Yan Ding⁸, David J. Dooling⁵, Shannon P. Dugan⁸, Wm Michael Dunne^{35,36}, A. Scott Durkin⁶, Robert C. Edgar³⁷, Rachel L. Erlich², Candace N. Farmer⁵, Ruth M. Farrell³⁸, Karoline Faust^{39,40}, Michael Feldgarden², Victor M. Felix⁷, Sheila Fisher⁸, Anthony A. Fodor⁴¹, Larry J. Forney⁴², Leslie Foster⁶, Valentina Di Francesco¹⁹, Jonathan Friedman⁴³, Dennis C. Friedrich², Catrina C. Fronick², Lucinda L. Fulton⁵, Hongyu Gao⁵, Nathalia Garcia⁴⁴, Georgia Giannoukos², Christina Gublin¹⁹, Maria Y. Giovanni¹⁹, Jonathan M. Goldberg², Johannes Goll⁶, Antonio Gonzalez⁴⁵, Allison Griggs², Sharvari Gujja², Susan Kienber Haake⁴⁶, Brian J. Haas², Holli A. Hamilton²⁹, Emily L. Harris²⁹, Theresa A. Herburn², Brandi Herter⁵, Diane E. Hoffmann⁴⁷, Michael E. Holder⁸, Clinton Howarth², Katherine H. Huang², Susan M. Huse⁴⁸, Jacques Izard^{32,33}, Janet K. Jansson⁴⁹, Huaiyang Jiang⁸, Catherine Jordan⁷, Vandita Joshi⁸, James A. Katancik⁵⁰, Wendy A. Keitel¹⁶, Scott T. Kelley⁵¹, Cristyn Kells², Nicholas B. King⁵², Dan Knights⁴⁵, Heidi H. Kong⁵³, Omry Koren⁵⁴, Sergey Koren⁵⁵, Karthik C. Kota⁵, Christie L. Kovar^{32,56}, Nikos C. Kyrpidis²⁷, Patricio S. La Rosa³⁴, Sandra L. Lee⁸, Katherine P. Lemon^{32,56}, Niall Lennon², Cecil M. Lewis⁵⁷, Lora Lewis⁸, Ruth E. Ley⁵⁴, Kelvin Li⁶, Konstantinos Liolios²⁷, Bo Liu⁵⁵, Yue Liu⁸, Chien-Chi Lo²⁶, Catherine A. Lozupone³, R. Dwayne Lunsford²⁹, Tessa Madden⁵⁸, Anup A. Mahurkar⁷, Peter J. Mannon⁵⁹, Elaine R. Mardis⁵, Victor M. Markowitz^{27,28}, Konstantinos Mavromatis²⁷, Jamison M. McCarrison⁶, Daniel McDonald³, Jean McEwen²¹, Amy L. McGuire⁶⁰, Pamela McInnes²⁹, Teena Mehta², Kathie A. Mihindukulasuriya⁵, Jason R. Miller⁷, Patrick J. Minx⁵, Irene Newsham⁶, Chad Nusbaum², Michelle O'Laughlin⁵, Joshua Orvis⁷, Ioanna Pagan²⁷, Krishna Palaniappan²⁸, Shital M. Patel⁶¹, Matthew Pearson², Jane Peterson²¹, Mircea Podar⁶², Craig Pohl³, Katherine S. Pollard^{63,64,65}, Mihai Pop^{55,66}, Margaret E. Priest⁷, Lita M. Proctor²¹, Xiang Qin⁸, Jeroen Raes^{39,40}, Jacques Ravel⁷, Jeffrey G. Reid⁸, Mina Rho⁶⁷, Rosamond Rhodes⁶⁸, Kevin P. Riehle⁶⁹, Maria C. Rivera^{23,24}, Beltran Rodriguez-Mueller³¹, Yu-Hui Rogers⁵, Matthew C. Ross¹⁶, Carsten Russ², Ravi K. Sanka⁶, Pamela Sankar⁷⁰, J. Fah Sathirapongsasuti¹, Jeffery A. Schloss²¹, Patrick D. Schloss⁷¹, Thomas M. Schmidt⁷², Matthew Scholz²⁶, Lynn Schriml⁷, Alyxandria M. Schubert⁷³, Nicola Segata¹, Julia A. Segre²⁹, William D. Shannon³⁴, Richard R. Sharp³⁸, Thomas J. Sharpton⁶³, Narmada Shenoy², Nihar U. Sheth²³, Gina A. Simone⁷³, Indresh Singh⁶, Christopher S. Smillie⁴³, Jack D. Sobel⁷⁴, Daniel D. Sommer⁵⁵, Paul Spicer⁵⁷, Granger G. Sutton⁶, Sean M. Sykes⁶, Diana G. Tabbaa², Mathangi Thiagarajan⁶, Chad M. Tomlinson⁵, Manolito Torralba⁶, Todd J. Treangen⁷⁵, Rebecca M. Truty⁶³, Tatiana A. Vishnivetskaya⁶², Jason Walker⁵, Lu Wang²¹, Zhengyuan Wang⁵, Doyle V. Ward², Wesley Warren⁵, Mark A. Watson³⁵, Christopher Wellington²¹, Kris A. Wetterstrand²¹, James R. White⁷, Katarzyna Wilczek-Boney⁸, Yuanqing Wu⁸, Kristine M. Wylie⁵, Todd Wylie⁵, Chandri Yandava², Liang Ye⁵, Yuzhen Ye⁶⁷, Shibu Yooshef⁷⁶, Bonnie P. Youmans¹⁶, Lan Zhang⁵, Yanjiao Zhou⁵, Yiming Zhu⁸, Laurie Zoloth⁷⁷, Jeremy D. Zucker², Bruce W. Birren², Richard A. Gibbs⁸, Sarah K. Highlander^{8,16}, Barbara A. Methé⁶, Karen E. Nelson⁶, Joseph F. Petrosino^{8,78,16}, George M. Weinstock⁵, Richard K. Wilson⁵ & Owen White⁷

¹Biostatistics, Harvard School of Public Health, Boston, Massachusetts 02115, USA. ²The Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02142, USA.

³Department of Chemistry and Biochemistry, University of Colorado, Boulder, Colorado 80309, USA. ⁴Howard Hughes Medical Institute, Boulder, Colorado 80309, USA. ⁵The Genome Institute, Washington University School of Medicine, St. Louis, Missouri 63108, USA. ⁶J. Craig Venter Institute, Rockville, Maryland 20850, USA. ⁷Institute for Genome Sciences, University of Maryland School of Medicine, Baltimore, Maryland 21201, USA.

⁸Human Genome Sequencing Center, Baylor College of Medicine, Houston, Texas 77030, USA. ⁹Department of Pathology & Immunology, Baylor College of Medicine, Houston, Texas 77030, USA. ¹⁰Department of Pathology, Texas Children's Hospital, Houston, Texas 77030, USA. ¹¹Department of Obstetrics & Gynecology, Division of Maternal-Fetal Medicine, Baylor College of Medicine, Houston, Texas 77030, USA. ¹²Molecular and Cellular Biology, University of Guelph, Guelph, Ontario N1G 2W1, Canada. ¹³Department of Civil & Environmental Engineering, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA. ¹⁴Center for Environmental Biotechnology, Lawrence Berkeley National Laboratory, Berkeley, California 94720, USA. ¹⁵School of Dentistry, University of California, San Francisco, San Francisco, California 94143, USA. ¹⁶Molecular Virology and Microbiology, Baylor College of Medicine, Houston, Texas 77030, USA.

¹⁷National Institute of Arthritis and Musculoskeletal and Skin, National Institutes of Health, Bethesda, Maryland 20892, USA. ¹⁸Office of Research on Women's Health, National Institutes of Health, Bethesda, Maryland 20892, USA. ¹⁹National Institute of Allergy and Infectious Diseases, National Institutes of Health, Bethesda, Maryland 20892, USA. ²⁰Department of Medicine, New York University Langone Medical Center, New York, New York 10016, USA. ²¹National Human Genome Research Institute, National Institutes of Health, Bethesda, Maryland 20892, USA. ²²Department of Statistical Sciences and Operations Research, Virginia Commonwealth University, Richmond, Virginia 23284, USA. ²³Center for the Study of Biological Complexity, Virginia Commonwealth University, Richmond, Virginia 23284, USA. ²⁴Department of Biology, Virginia Commonwealth University, Richmond, Virginia 23284, USA. ²⁵Technology Integration Group, National Energy Research Scientific Computing Center, Lawrence Berkeley National Laboratory, Berkeley, California 94720, USA. ²⁶Genome Science Group, Bioscience Division, Los Alamos National Laboratory, Los Alamos, New Mexico 87545, USA. ²⁷Joint Genome Institute, Walnut Creek, California 94598, USA. ²⁸Biological Data Management and Technology Center, Computational Research Division, Lawrence Berkeley National Laboratory, Berkeley, California 94720, USA. ²⁹National Institute of Dental and Craniofacial Research (NIDCR), National Institutes of Health, Bethesda, Maryland 20892, USA. ³⁰FemCare Product Safety and Regulatory Affairs, The Procter & Gamble Company, Cincinnati, Ohio 45224, USA. ³¹Bioinformatics Department, Second Genome, Inc., San Bruno, California 94066, USA. ³²Department of Molecular Genetics, Forsyth Institute, Cambridge, Massachusetts 02142, USA. ³³Department of Oral Medicine, Infection and Immunity, Harvard School of Dental Medicine, Boston, Massachusetts 02115, USA. ³⁴Department of Medicine, Division of General Medical Science, Washington University School of Medicine, St. Louis, Missouri 63110, USA. ³⁵Department of Pathology & Immunology, Washington University School of Medicine, St. Louis, Missouri 63110, USA. ³⁶bioMerieux, Inc., Durham, South Carolina 27712, USA. ³⁷drive5.com, Tiburon, California 94920, USA. ³⁸Center for Ethics, Humanities and Spiritual Care, Cleveland Clinic, Cleveland, Ohio 44195, USA. ³⁹Department of Structural Biology, VIB, Belgium, 1050 Ixelles, Belgium. ⁴⁰Department of Applied Biological Sciences (DBIT), Vrije Universiteit Brussel, 1050 Ixelles, Belgium. ⁴¹Department of Bioinformatics and Genomics, University of North Carolina - Charlotte, Charlotte, North Carolina 28223, USA. ⁴²Department of Biological Sciences, University of Idaho, Moscow, Idaho 83844, USA. ⁴³Computational and Systems Biology, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA. ⁴⁴Center for Advanced Dental Education, Saint Louis University, St. Louis, Missouri 63104, USA. ⁴⁵Department of Computer Science, University of Colorado, Boulder, Colorado 80309, USA. ⁴⁶Division of Associated Clinical Specialties and Dental Research Institute, UCLA School of Dentistry, Los Angeles, California 90095, USA. ⁴⁷University of Maryland Francis King Carey School of Law, Baltimore, Maryland 21201, USA. ⁴⁸Josephine Bay Paul Center, Marine Biological Laboratory, Woods Hole, Massachusetts 02543, USA. ⁴⁹Ecology Department, Earth Sciences Division, Lawrence Berkeley National Laboratory, Berkeley, California 94720, USA. ⁵⁰Department of Periodontics, University of Texas Health Science Center School of Dentistry, Houston, Texas 77030, USA. ⁵¹Department of Biology, San Diego State University, San Diego, California 92182, USA. ⁵²Faculty of Medicine, McGill University, 3647 Peel St, Montreal, Quebec H3A 1X1, Canada. ⁵³Dermatology Branch, CCR, National Cancer Institute, Bethesda, Maryland 20892, USA. ⁵⁴Department of Microbiology, Cornell University, Ithaca, New York 14853, USA. ⁵⁵Center for Bioinformatics and Computational Biology, University of Maryland, College Park, Maryland 20742, USA. ⁵⁶Division of Infectious Diseases, Children's Hospital Boston, Harvard Medical School, Boston, Massachusetts 02115, USA. ⁵⁷Department of Anthropology, University of Oklahoma, Norman, Oklahoma 73019, USA. ⁵⁸Department of Obstetrics and Gynecology, Washington University School of Medicine, Saint Louis, Missouri 63110, USA. ⁵⁹Division of Gastroenterology and Hepatology, University of Alabama at Birmingham, Birmingham, Alabama 35294, USA. ⁶⁰Center for Medical Ethics and Health Policy, Baylor College of Medicine, Houston, Texas 77030, USA. ⁶¹Medicine-Infectious Disease, Baylor College of Medicine, Houston, Texas 77030, USA. ⁶²Biosciences Division, Oak Ridge National Laboratory, Oak Ridge, Tennessee 37831, USA. ⁶³Gladstone Institutes, University of California, San Francisco, San Francisco, California 94158, USA. ⁶⁴Institute for Human Genetics, University of California, San Francisco, San Francisco, California 94158, USA. ⁶⁵Division of Biostatistics, University of California, San Francisco, San Francisco, California 94158, USA. ⁶⁶Department of Computer Science, University of Maryland, College Park, Maryland 20742, USA. ⁶⁷School of Informatics and Computing, Indiana University, Bloomington, Indiana 47405, USA. ⁶⁸Mount Sinai School of Medicine, New York, New York 10029, USA. ⁶⁹Molecular & Human Genetics, Baylor College of Medicine, Houston, Texas 77030, USA. ⁷⁰Center for Bioethics and Department of Medical Ethics, University of Pennsylvania, Philadelphia, Pennsylvania 19104, USA. ⁷¹Department of Microbiology & Immunology, University of Michigan, Ann Arbor, Michigan 48109, USA. ⁷²Department of Microbiology and Molecular Genetics, Michigan State University, East Lansing, Michigan 48824, USA. ⁷³The EMMES Corporation, Rockville, Maryland 20850, USA. ⁷⁴Harper University Hospital, Wayne State University School of Medicine, Detroit, Michigan 48201, USA. ⁷⁵McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University School of Medicine, Baltimore, Maryland 21205, USA. ⁷⁶J. Craig Venter Institute, San Diego, California 92121, USA. ⁷⁷Feinberg School of Medicine, Northwestern University, Chicago, Illinois 60611, USA. ⁷⁸Alkek Center for Metagenomics and Microbiome Research, Baylor College of Medicine, Houston, Texas 77030, USA. ⁷⁹Genetics and Molecular Biology Branch, National Human Genome Research Institute, Bethesda, Maryland 20892, USA.

*These authors contributed equally to this work.